

A Neural Network Model of Episode Representations in Working Memory

Martin Takac^{1,2} · Alistair Knott¹

Received: 2 September 2014 / Accepted: 20 April 2015
© Springer Science+Business Media New York 2015

Abstract We present a neural network model of the storage of episode representations in working memory (WM). Our key idea is that episodes are encoded in WM as prepared sensorimotor routines, i.e. as prepared sequences of attentional and motor operations. Our network reproduces several experimental findings about the representation of prepared sequences in prefrontal cortex. Interpreted as a model of WM episode representations, it has useful applications in an account of long-term memory for episodes and in accounts of sentence processing.

Keywords Working memory · Neural network modelling · Sequence learning · Action preparation · Language processing

Introduction

Working Memory for Episodes

The classical model of working memory (WM) by Baddeley and Hitch [1] posits two representational media: one for visual material (the visuospatial sketchpad) and one for phonological material (the phonological buffer). Baddeley [2] revised the model to include a new third medium,

holding semantic material, called the ‘episodic buffer’. This medium stores semantic representations of actions, or events, or stative propositions: we will use the term *episodes* to refer to these representations. Our paper is about the episodic buffer.

Baddeley argues for the episodic buffer on several grounds. One argument concerns the process through which episode representations are stored in long-term memory. The long-term neural storage of an episode is widely agreed to involve the hippocampus [3]: specifically, the creation of links between hippocampal assemblies representing the various semantic components of the episode.

But associations between hippocampal assemblies can only be learned if they are active in quick succession, within around 100 ms of one another [4]. Experiencing an episode often takes much longer than this. So, we must envisage that episode representations are initially buffered in some WM medium and are only relayed from there to the hippocampus when they are complete.

Another argument for the episodic buffer relates to language processing. When a speaker produces a sentence, the message to be expressed is standardly assumed to be held in a WM medium holding semantic representations (see, for example, Levelt [5]). When a hearer interprets a sentence, the inferred meaning is likewise assumed to be assembled in a WM medium holding semantic representations (e.g. [6, 7]). Baddeley [2] suggests that the episodic buffer is also the medium holding semantic representations in sentence processing tasks; accordingly, he suggests that it is connected bidirectionally to the phonological buffer, so that words and semantic representations can activate one another during sentence processing. Thus, the episodic buffer is accorded a role in language processing as well as a role in creating long-term memories.

✉ Martin Takac
takac@ii.fmph.uniba.sk

Alistair Knott
alick@cs.otago.ac.nz

¹ Department of Computer Science, University of Otago,
PO Box 56, 9054 Dunedin, New Zealand

² Centre for Cognitive Science FMFI UK, Comenius
University, Mlynska Dolina, 842 48 Bratislava, Slovakia

In this paper, we present a connectionist model of WM storage which supports both memory and language-processing functions of the episodic buffer. Our main novel proposal is that episodes are stored in WM as *prepared sequences of attentional and motor operations* (see the next section). This proposal places WM episodes within a class of mental representation that is relatively well understood: we know a lot about how the primate brain (in particular, the prefrontal cortex) stores prepared sequences of attentional and motor movements, from single-cell studies in monkeys. The proposal also gives a natural account of how WM episode representations are communicated to the hippocampus and how they interface with language: if episodes are represented as prepared sequences, they can be transmitted to other cognitive modalities by being *replayed*.

In the next two sections, we introduce our proposal that episodes are stored in WM as prepared sensorimotor sequences and we review evidence from monkey studies that prepared sensorimotor sequences are stored in prefrontal cortex, in a number of distinct formats.

WM Episode Representations as Prepared Sensorimotor Routines

Our model is founded on the assumption that WM episodes provide an interface between the sensorimotor mechanisms through which episodes are apprehended and the hippocampal structures in which they are stored. On this assumption, we expect the structure of WM episode representations to reflect both the structure of sensorimotor processes and the structure of hippocampal representations. A strong commonality in the structures of these two domains is *sequential organisation*.

Sensorimotor processing is strongly sequential at certain timescales, because it involves sequential deployments of the agent's sensory and motor apparatus. (For instance, saccades deploy the agent's fovea sequentially to targets in the world.) Ballard et al. [8] propose that sensorimotor processing is organised into sequentially structured routines, whose atomic elements are discrete sensory or motor actions. These actions are termed *deictic operations*, and a sequence of such actions is termed a *deictic routine*. Through a case study of episodes involving reach-to-grasp actions, Knott [9, 10] argues that the sensorimotor processes through which concrete episodes are apprehended take the form of sequentially structured deictic routines.

The hippocampus stores associations between stimuli of many different kinds. But an emerging idea is that it is specially good at storing associations between sequentially structured items [11]. One finding which strongly supports this idea is that the hippocampus actively *replays* sequences of representations evoked during sensorimotor experience [12] (see [13] for a model). The key result is

that sequences of hippocampal place cells activated when a rat navigates a maze are replayed later when the rat is asleep. (Sequences are replayed at much higher speeds, consistent with the hippocampus' natural recurrent dynamics.) Since episodes appear to be apprehended through well-defined sequences of sensorimotor operations and sequences appear to be a natural unit of storage in the hippocampus, an interesting possibility is that WM episodes are also stored as sequences. Our model of WM episodes basically implements this idea.

Our core proposal is that a concrete episode is stored in WM as a representation of the sequence of sensorimotor operations through which it was experienced. This proposal suggests a novel style of solution to a well-known problem for models of neural episode representation. An episode representation must include representations of the individual participants in an episode, but must also *bind* each participant representation to the particular role it plays in the episode (e.g. 'agent' or 'patient'). There is no consensus as to how this binding is implemented in the brain. (We will review some of the suggestions in "[Comparison with Other Models of Episode Representation](#)".) Our novel proposal is that the order of sensorimotor operations in a deictic routine implicitly identifies the roles played by participants in the observed episode, because participants playing particular roles are attended to at canonical positions in the routine.

This position is argued for in detail by Knott in [9]. We will illustrate by briefly considering the process involved in experiencing a man grabbing a cup. If the observer is the man performing the action, the first operation he must perform is to decide to act [14]; it is only after this decision is made that his motor system is configured for action execution [15]. The decision to act activates a representation of the self as the agent of the forthcoming action [14, 16]. Having decided to act, the agent cannot immediately execute a transitive motor action; a target for this action must first be selected [17] and then attended to, so that its detailed motor affordances can be computed [18]. Only then can a specific motor programme be selected and executed. On this analysis, experiencing a reach-to-grasp episode from the perspective of the agent involves a strict sequence of sensorimotor operations: ATTEND-TO-AGENT, ATTEND-TO-TARGET, ACTIVATE-GRASP-PROGRAMME. A similar analysis is given for the process involved in perceiving an external agent performing a grasp action; see Knott [9] for details. There is also evidence for distinct sequences of sensorimotor operations associated with experiencing communicative actions [19, 20], causative actions such as opening and breaking [21], and actions involving moving oneself or a controlled object along a trajectory.¹

¹ Manuscript in preparation, Knott and Takac: Locomotion actions as sequentially structured sensorimotor routines.

Our specific proposal is that episodes are stored in WM as *prepared sequences of attentional and motor operations*, in a medium in prefrontal cortex that holds prepared behavioural sequences. This proposal is roughly in line with Baddeley's own sketched idea that the visuospatial sketchpad and episodic buffer are rehearsed by a general process of 'sequential attention' (Baddeley [2: 420]), and with Curtis and D'Esposito's [22] general proposal that WM representations are held in high-level motor control regions. But more importantly, it has several advantages in its own right. For one thing, as already mentioned, we know something about *how* prepared sequences are stored in this prefrontal medium, so the proposal places WM episode representations within a class of neural representation about which something is already known. In addition, the proposal fits well with simulationist theories of meaning, of the kind that feature prominently in embodied models of cognition. In these theories, activating a semantic episode representation involves actively simulating the process of experiencing the episode (see, for example, Barsalou [23] or Glenberg and Gallese [24]). If a WM episode is a prepared sensorimotor routine, it is naturally something that can be executed, either in overt behaviour or in simulation, so there is a natural implementation of the process of simulating an episode. Finally, our proposal allows a natural account of how WM episodes are transmitted to long-term memory and are expressed linguistically. If WM episodes can be replayed, we can envisage that transmission involves a replay operation, taking place in a mode of brain connectivity where sensorimotor signals activate assemblies in some selected cognitive medium, at a speed commensurate with that medium. An episode stored as a prepared sequence in WM can be transmitted to the hippocampus by being replayed fast, at around the 100 ms timescale, with the hippocampus configured to encode an incoming sequence of stimuli; or it can be relayed to the speech production system by being replayed at a timescale commensurate with the production of a stream of words, in a mode where sensorimotor signals activate output phonology. In summary, in our proposal, episodes are experienced as sequences and stored in WM as prepared sequences; they can then be replayed to the hippocampus where they are stored more permanently as sequences, or replayed to the linguistic system, where they generate sequences of words.

Since we envisage WM episode representations as having a role in the creation of longer-term episode representations in the hippocampus, it is useful to situate our model in relation to the classic account of hippocampal episodic memory (see, for example, Tulving [25]). For one thing, while Tulving's model concerns the storage of episodes in LTM, our model focusses on the storage of episodes in WM, *prior* to their storage in the hippocampus. For

another thing, while Tulving's model emphasises the sequential structure of episodic memory, the items organised into sequences in his model are *whole episodes*; our model focusses on sequential structures *within single episodes*.

Representation of Prepared Sequences in Prefrontal Cortex

A bonus of the model just outlined is that the neural mechanisms supporting preparation of sensorimotor sequences have been extensively studied, in single-cell recording experiments in monkeys. The principal mechanisms supporting sequence preparation are in dorsolateral prefrontal cortex (dlPFC) [26, 27]. Several schemes for encoding prepared sequences have been found. In this section, we will review these.

In one scheme, individual neurons encode specific movements in particular contexts. For instance, Barone and Joseph [26] found neurons which were active when a monkey prepared movement *A*, but only when it was followed by another movement *B*. We could call these representations 'place-coded' representations of sequences.

In another scheme, neurons encode individual movements, and their position in the prepared sequence is given by their activation levels. For instance, in a monkey preparing a sequence of three movements *A*, *B*, and *C*, Averbeck et al. [27] found neurons representing each prepared action which were active in parallel, with the neuron encoding *A* most active and that encoding *C* least active. Interestingly, when the prepared sequence is executed, neurons encoding specific actions using this scheme are inhibited just after their associated action is produced. Averbeck et al.'s [27] findings strongly support a 'competitive queueing' model of sequence preparation, in which PFC assemblies encoding different actions compete against one another, with the winner triggering the associated action, but also an operation to inhibit itself, so the next-most active assembly wins the competition at the next time point [28]. In competitive queueing, the representation of a prepared sequence is destructively updated in the medium in which competition occurs. We will call the sequence representations in this medium 'dynamic'. However, there is also evidence that prepared sequences are represented in a WM medium which is *not* destructively updated when a sequence is replayed. A simple argument for this is that a given prepared sequence can be executed several times: each time, the sequence representation in the dynamic medium must somehow be restored from some more enduring medium. We will call representations in the enduring medium 'static'.

There is also evidence that a monkey can represent multiple alternative prepared sequences in dlPFC, in a medium which allows competition between candidate sequences and the selection of a winner. This evidence comes

from a study of Averbeck et al. [29], in which monkeys were trained to perform two sequences in response to two cues. Each day, different cues were chosen to represent the two sequences. Halfway through the day, the mapping from cues to sequences was reversed, so the monkeys had to gradually learn the new mapping. During this period, dlPFC assemblies could be identified representing each prepared sequence, and the relative activation of the two assemblies after presentation of a cue could be used to predict the sequence which the monkey actually performed.

In summary, the prefrontal mechanism implementing sequence preparation appears to involve four distinct media. There is a medium holding representations of individual operations in a sequence, which encodes the context in which they appear. There is a medium holding distributed representations of whole sequences, in assemblies whose components encode individual actions, whose order is determined by their level of activation. Sequence representations in this medium are destructively updated when a prepared sequence is executed. But there is also a medium holding sequence representations which are not destroyed. Finally, there is a medium in which alternative candidate sequence representations are active in parallel and compete with one another. If episodes are stored in WM as prepared sensorimotor sequences, then this mechanism would allow for WM episodes to be stored and replayed and also for alternative WM episodes to compete amongst one another, with the winner being selected.

There is also good evidence that dorsolateral PFC is involved in planning sequences of actions in humans. This area is active when subjects maintain complex prepared actions in WM (see, for example, [22]), and damage or inactivation of this area results in impairments in planning performance (see, for example, [30, 31]). However, this evidence is much less detailed than the single-neuron studies, so we will rely mainly on the macaque studies in developing our computational model.

Methods

A Neural Network Model for Storage and Selection of WM Episodes

In this section, we describe a neural network model of the sequence preparation mechanisms outlined in the previous section.

We envisage the network being used to control the process of ‘experiencing an episode’ both when the experiencer is acting himself, and thus ‘bringing about’ the episode, and when he is passively watching an external episode. In the former case, the mechanisms that ‘select’ an episode should be thought of as planning mechanisms that

decide what action to do; in the latter case, the same selection mechanisms should be thought of as generating expectations about the episodes that will happen next. Thus, the function of the episode selection mechanism will depend on the mode in which the network is being used. In presenting the network, we will focus on cases where it is used during perception.²

Our key aim for the network is that it learns the kind of representations of prepared sequences which are found in monkey PFC, as discussed above. However, there are also two other design criteria. Firstly, we want there to be a medium in which candidate sensorimotor operations compete with one another at every stage during the execution of a sequence. At any point, the operation which an agent executes is dictated partly by what is planned or expected, but also partly by bottom-up stimuli. We want a medium which allows competition between alternative operations from both these sources. Secondly, it is important that the medium representing alternative possible sensorimotor sequences avoids binding errors, whereby an operation belonging to one sequence is falsely identified as part of a different sequence. Given that this medium must represent multiple sequences simultaneously, this is a difficult requirement. To address both these criteria, a key design decision is to use self-organising maps (SOMs) [33], both to represent the individual sensorimotor operations presented to the system and to represent sequences of these operations, as explained in the following section.

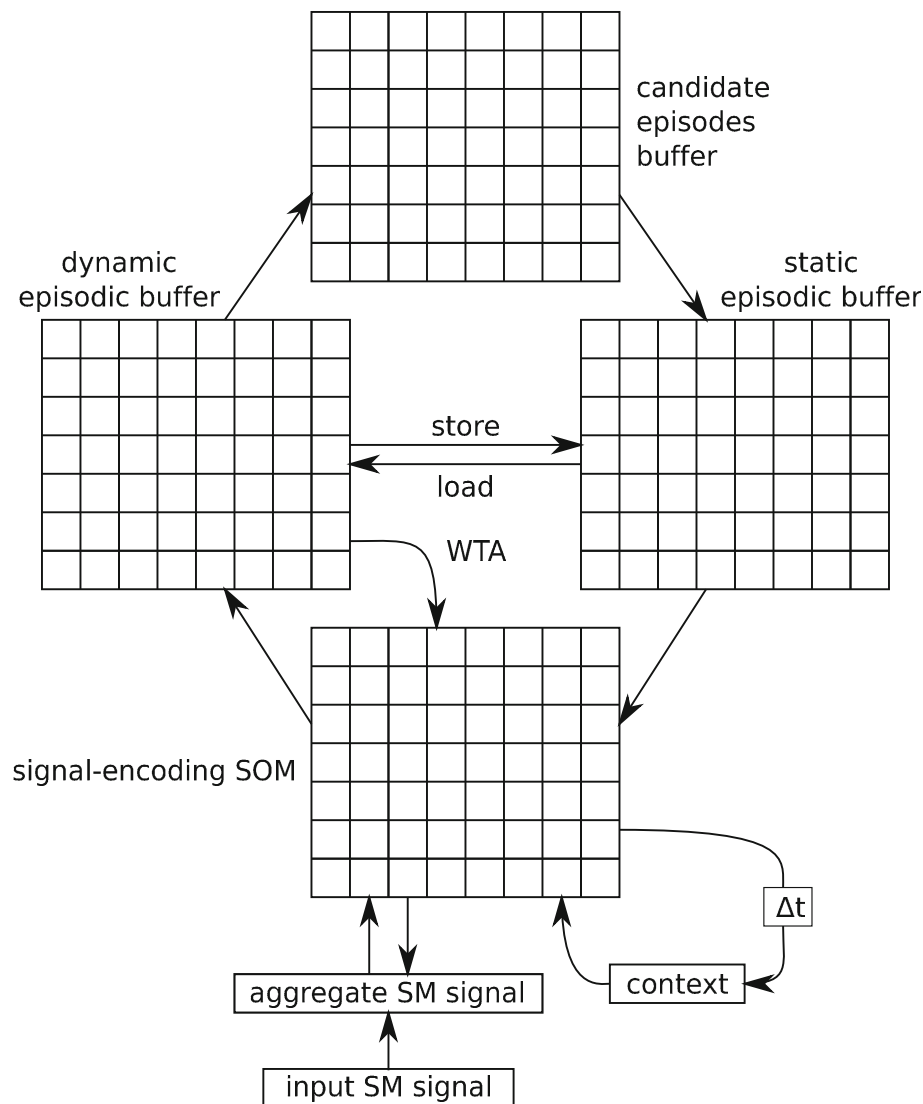
Architecture

The architecture of our network is shown in Fig. 1. The network takes as input a sequence of sensorimotor signals at successive time points, evoked in the *input SM signal* area. Input SM signals can be motor actions (representing actions of the agent or of a perceived external agent) or attentional actions (representing objects).³ In most of the

² In fact, even when the network is used for perception of episodes, the selection mechanism has an impact on the agent’s behaviour. The expected episode is a sequence whose first two items are planned *attentional actions*: during perception, these actions will actually be executed by the agent, with results that depend on the world as well as on the agent’s expectations and which might well result in revisions to the selected episode. The model thus allows for active perceptual operations during the process of selecting an episode, creating the structural coupling between its representational system and the environment that is characteristic of embodied systems (see, for example, [32]).

³ The most obvious kind of ‘attentional actions’ are overt movements, such as saccades. But attentional actions also involve cognitive operations, in particular top–down activation of semantic representations. These top–down activation operations can encode either the expected result of a forthcoming object classification process [34] or the category of properties of a desired search target [35, 36]. Our attentional actions represent objects in the sense that they represent expected or sought-for object categories.

Fig. 1 Architecture of the network. A temporal sequence of input SM signals causes an activity in the recurrent signal-encoding SOM. The winning units leave an exponentially decaying activity trace in the dynamic episodic buffer via excitatory connections. After the sequence is completed, the trace is stored in the static episodic buffer via excitatory connections between isomorphic units. The trace also serves as training input for the candidate episodes buffer, which can provide top-down bias (expectations) via static episodic buffer back to the dynamic one. The trace also serves as training input for the candidate episodes buffer, which can provide top-down bias (expectations) via static episodic buffer back to the dynamic one



experiments in this paper, we model individual SM inputs as localist units. However, we also conducted some initial experiments with distributed inputs, which are described in a separate technical report [37]. In both cases, the experience of an episode produces a sequence of well-delineated, discrete SM signals in the input area. Obviously, these input stimuli greatly simplify the actual sensorimotor stimuli that arise during experience of an episode, particularly in the localist implementation. However, while we are undoubtedly simplifying, we argue that the clearly structured *temporal* organisation of the inputs is a genuine feature of SM processing: as already discussed, our contention is that deictic routines deliver a reasonably discrete sequence of SM representations. Our model is ‘embodied’ in that it leverages this temporal organisation, rather than the internal composition of individual SM signals.

Input SM signals are fed through an *aggregate SM signal* area (described below) to a *signal-encoding SOM*.

This SOM has recurrent connections, as described by Strickert and Hammer [38]: it takes as an additional input a set of context units, whose activations reflect the weight vector of the winning unit at the previous time point. (The weight vector of each unit now includes weights of connections from the context units, which is what creates recurrency in the network.)

When trained on a sequence of inputs, a recurrent SOM organises itself so that individual units encode signals occurring in particular sequential contexts, very much like the PFC units identified by Barone and Joseph [26]. This learning happens gradually over time and is implemented in synaptic strengths rather than activations, so the signal-encoding SOM does not hold WM episode representations itself—but it provides the inputs from which WM episode representations are created.

Units in the signal-encoding SOM represent signals in a localist way, so that alternative signals compete with one

another. The winning signal at each time step is copied to an area which is isomorphic with the recurrent SOM called the *dynamic episodic buffer* (see Fig. 2). This area accumulates representations of each signal in an input sequence, with the first signal represented most strongly and subsequent signals being stored with decreasing activation, as in the prefrontal area studied in [27]. The dynamic episodic buffer holds an activity-based representation of a single episode in WM. When an input sequence is encoded in the dynamic episodic buffer, it can be replayed immediately by iteratively sending the dynamic episodic buffer's most active unit to the signal-encoding SOM (via the 'WTA' link) and then inhibiting this winning unit (a process we loosely term 'inhibition-of-return'). To support repeated execution of a sequence, it can be stored in a *static episodic buffer*, which has the same structure as the dynamic one, and later reloaded. The static episodic buffer is also an activity-based medium: it holds a complete episode representation in WM as an active pattern.

At the highest level in the network, there is another SOM called the *candidate episodes buffer*. This area takes as its input the distributed representations in the dynamic episodic buffer and encodes them as localist units. During training, it learns to represent episodes with similar encodings in the dynamic episodic buffer in neighbouring positions in the SOM. At every time point during presentation of a sequence, this area represents a probability distribution over complete episodes. (If the network is being used to control the agent's own actions, this distribution represents action sequences which lead to reward; if it is being used to support observation of external episodes, it represents likely action sequences.) The distribution changes as new items arrive in the sequence and become encoded in the dynamic episodic buffer. Like the signal-encoding SOM, the candidate episodes buffer does not hold WM episode representations itself; rather it provides a learned top-down bias that influences which episode is held in WM at any given time.

The winning unit in the candidate episodes buffer provides top-down activation to the static episodic buffer, through weights which are copies of those delivering input to the candidate episodes buffer. Since the winning unit always encodes a complete episode, the static episodic buffer likewise always encodes a complete episode, but in the same distributed format used by the dynamic episodic buffer. During presentation of a sequence, activity in the static episodic buffer is fed back to the signal-encoding SOM. This top-down input, when combined with the current context representation, produces a pattern of activity biased towards a representation of the next SM signal. The pattern is passed back to the aggregate SM signal area at the next time point. Thus, the aggregate area receives both bottom-up inputs from the input SM signal and top-down ones from the static episodic buffer. We think of the

aggregate area as the medium in which top-down goals or expectations modulate bottom-up perceptual inputs, but in the current model, this modulation is not implemented: during training, all inputs to the area are bottom-up, and during testing, all inputs are top-down.

Once a winner is selected in the candidate episodes buffer, activity is propagated back through the network, a process we call 'top-down reconstruction'. This process makes use of the special property of SOMs that the weights of each unit explicitly encode the input pattern it is most responsive to, allowing this pattern to be recreated top-down. During reconstruction, the weights of the winning unit in the candidate episodes buffer are copied back to the static and then dynamic episodic buffer. Destructive iterative updating of the dynamic episodic buffer causes a temporal sequence of activations of units in the signal-encoding SOM, which in turn projects their weight vectors back to the aggregate SM signal layer where they represent top-down expectations.

Not all units in the candidate episodes buffer allow a meaningful episode to be reconstructed. In particular, units that have never been selected as a winner during training are unlikely to have weights that allow the reconstruction of a meaningful episode. If their weights have been trained at all, it is because they are in the neighbourhood of units that have been winners in the past. (They may be in the neighbourhood of several such units, in which case we term them 'blend' units.) We want to bias competition in the candidate episodes buffer towards units that have been selected as winners. In our implementation, each unit in the candidate episodes buffer is associated with a second (scalar) weight called the prior. During training, units compete the standard way as is usual in SOMs, i.e. only based on the distance of their weight vectors to the input. Each time a unit becomes a winner, its prior is increased (the update rule for original weights stays unchanged). During prediction, the distance-based competition between units is biased by their priors: the activity of each unit is a weighted combination of its inverse distance from the winning unit and its prior. Thus, units with higher priors have a higher chance of becoming the winner.

Now, we report some technical details of the network architecture. The *input SM signal* layer consists of 35 units; SM signals are encoded with 1-hot localist coding, i.e. there is one unit for each possible SM signal. The *aggregate SM signal layer* is isomorphic with the input layer.

The *signal-encoding SOM* is a two-dimensional Merge SOM [38] with 400 units and parameters $\alpha = 0.4$, $\beta = 0.5$, constant learning rate $\gamma = 0.1$, and Gaussian neighbourhood with linearly decreasing σ from 10 to 0.5 in 25,000 steps [see Eq. (4) below].

The *dynamic episodic buffer* and the *static episodic buffer* are both 2D maps (not SOMs) isomorphic with the

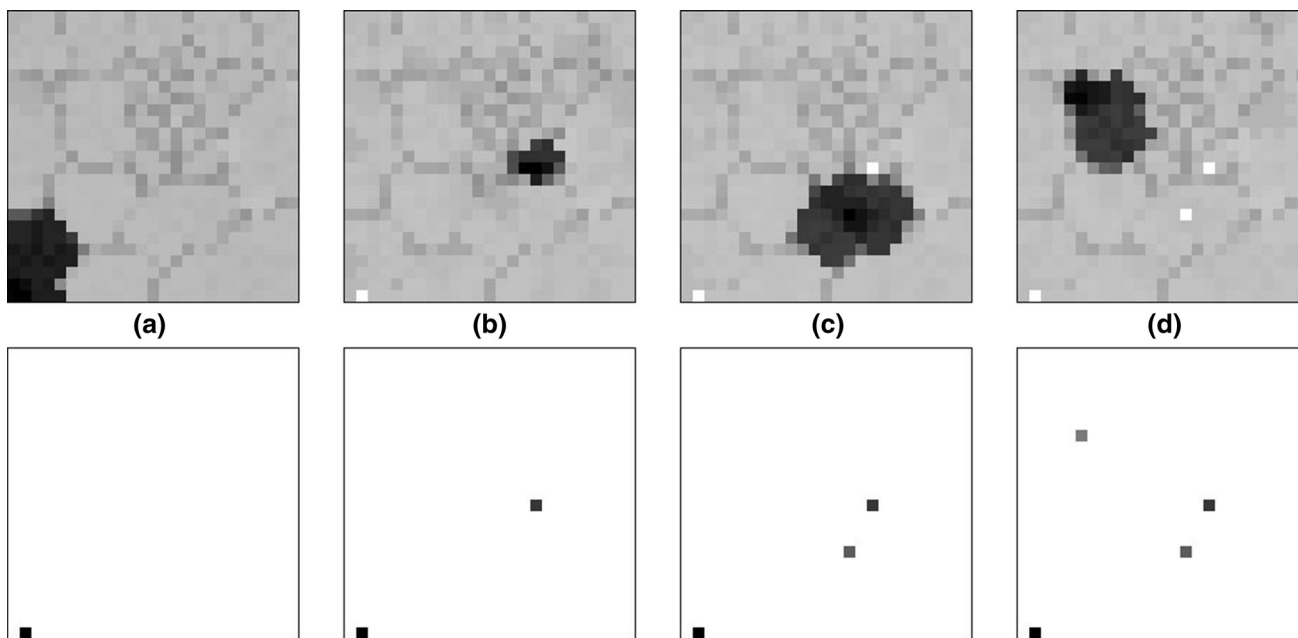


Fig. 2 Activity in the signal-encoding SOM (*top row*) and the trace it leaves in the dynamic episodic buffer (*bottom row*) during sequential presentation of an episode consisting of SM elements DOG **(a)** CHAIR

(b) SLEEP **(c)** UNDER **(d)**. The activities are *colour-coded* in shades of grey from white (activity = 0) to black (activity = 1)

signal-encoding SOM, i.e. they have 400 units each. Each unit in the dynamical episodic buffer is connected to the corresponding unit of the signal-encoding SOM.

Experiencing a sequence of SM operations creates a temporal pattern of active units in the signal-encoding SOM. Activity of the i -th unit at time t is inversely proportional to a combined squared Euclidean distance $\text{dist}_i(t)$ between its regular weight w_i and the input vector $x(t)$, and between the context weight c_i and the recursive context descriptor $c(t)$ (for details, see [38]):

$$\text{dist}_i(t) = (1 - \alpha) \cdot \|x(t) - w_i\|^2 + \alpha \cdot \|c(t) - c_i\|^2 \quad (1)$$

This activity is further decreased by the current activity $D_i(t)$ of the corresponding unit in the dynamical episodic buffer via an inhibitory connection, so that the resulting activity $A_i(t)$ of the i -th unit of the signal-encoding SOM is

$$A_i(t) = \max\left(0, e^{-\text{dist}_i(t)} - D_i(t)\right) \quad (2)$$

Then, the unit I with the highest activity A_I activates its isomorphic unit in the dynamical episodic buffer via an excitatory connection modulated by an independent exponentially decaying temporal signal δ^t , so that the sequence of winning units leaves a ‘trace’ of exponentially decaying activities in the isomorphic units in the dynamic episodic buffer:

$$D_i(t + 1) = \begin{cases} \delta^t & \text{if } i = I, \\ D_i(t) & \text{otherwise} \end{cases} \quad (3)$$

where $\delta = 0.8$ and all units start with zero activity $D_i(0) = 0$ (time t denotes the number of elements in the sequence seen so far). The inhibition from the dynamic episode buffer forces the signal-encoding SOM to select a new winner in each step of the sequence, in order to prevent the confusion of elements in the trace. The static episodic buffer is connected to the dynamic episodic buffer via excitatory connections between corresponding isomorphic units so that the activity from one buffer can be copied to the other one.

The weights of the signal-encoding SOM are updated using the standard SOM learning rule

$$w_i(t + 1) = w_i(t) + \gamma \cdot N(I, i) \cdot [x(t) - w_i(t)] \quad (4)$$

where N is a Gaussian neighbourhood function $N(I, i) = \exp(-\|r_I - r_i\|^2 / \sigma^2)$, I is the index of the winning neuron, and r_I, r_i are vectors of lattice coordinates of neurons I, i .

After completing the whole sequence, the 400-dimensional vector representing its trace serves as a training input to the *candidate episodes buffer* with 900 units, constant learning rate $\gamma = 0.9$, and Gaussian neighbourhood with linearly decreasing σ from 10 to 0.5 [see Eq. (4)]. The candidate episodes buffer is a standard SOM with modifications mentioned above: each unit maintains a scalar prior weight p_i that is increased each time this unit becomes the winner. The result of this is

$$p_i = 1 - e^{-c \cdot N} \quad (5)$$

where N is the number of times the unit became the winner and $c = 0.01$ is a scaling constant. During prediction, competition between active units is biased by prior weights in that the activity C_i of the i -th unit is computed as

$$C_i = (1 - \alpha) \cdot e^{-d_i} + \alpha \cdot p_i \quad (6)$$

where d_i is a standard SOM Euclidean distance between the unit's weight vector and the SOM's input and $\alpha = 0.35$ is a mixing coefficient.

Training

We trained the model on sequences of sensorimotor signals, representing the sensorimotor routines through which different episodes are experienced. The sequences were built from 35 sensorimotor signals corresponding to six attentional actions (representing objects MAN, DOG, CAT, CUP, BALL, and CHAIR), 25 motor programs (e.g. SIT, GRAB, PUSH, CAUSE, BREAK), three prepositional relations (UNDER, BEHIND, and NEAR), and an end-of-episode signal (.).⁴ The order of signals in a sequence was not arbitrary, but was determined by the type of a represented episode, e.g. MAN SNEEZE (intransitive episode), MAN CUP GRAB (transitive), MAN SLEEP CHAIR NEAR (intransitive with prepositional phrase complement), MAN CUP CAUSE BREAK (simple causative), and DOG BALL CAUSE GO CHAIR UNDER (causative with prepositional phrase). Detailed justification of the orderings in these sequences exceeds the scope of this paper, but it can be found in [9, 20, 21].

We repeated each simulation 10 times with different random initialisations of connection weights in the model and different training sets. Each training set consisted of 500 sequences stochastically generated by the same set of transcription rules (see Table 1). The training sets contained on average 2.6 % sequences of length 2, 17.3 % of length 3, 25.2 % of length 4, and 54.9 % of length 6. Sequences could contain duplicates: in all, 19.1 % of sequences contained two copies of a single signal and 0.9 % contained 3. The training took 200 epochs; in each epoch, the training sequences were presented in random order and the Merge SOM context was reset after each sequence.

⁴ Since WM representations are normally understood to be 'maintained in the face of incoming perceptual stimuli', we have to assume a special operation to remove the current episode representation from the dynamic episodic buffer before the next episode begins—an operation that probably involves an element of self-inhibition (see, for example, Mayr and Keele [39]). In previous work [19], we have considered the nature of this operation, but in the present study, we just use the end-of-episode signal to stand in for this operation.

Results

After training, we tested the network in two tasks: immediate serial recall and prediction of complete episodes from their fragments. All tests were repeated for the 10 different simulation runs and averaged.

Immediate Serial Recall

The basic requirement for our network is that it can store and replay individual behavioural sequences. This capability relies on interactions between the signal-encoding SOM and the dynamic episodic buffer. We presented the trained network with 200 sequences of input signals: 100 taken from the training data and 100 new ones not seen before. Each sequence was coded in the dynamic episodic buffer; then, the signal-encoding SOM's context was reset and the winning unit in the dynamic buffer was iteratively sent to the SOM and then inhibited. In total, 99.9 % (SD = 0.3 %) of training sequences were correctly replayed, and 98.9 % (SD = 1.14 %) of unseen sequences.

Note that the network is not confused by sequences containing duplicate items. A regular competitive queuing network has problems representing duplicate items, because after the first instance of the item is presented, it is inhibited in the competitive medium. But since the dynamic episodic buffer receives inputs from the signal-encoding SOM where we forced a unique winner selection, different instances of a given input are represented differently, and it does not suffer from this problem. To verify this, we also tested the immediate recall on a set consisting of 95 sequences with two repeating elements and five sequences with three repeating elements. The network correctly reconstructed all the sequences.

Predicted Completions of Sequences

The network is also designed to generate top-down predictions about sequences being experienced, through activity in the candidate episodes buffer. The prediction is actually a retrieval of a past episode most similar to the episode fragment experienced so far. Past episodes are remembered in the weights of units of the candidate episodes buffer. The fragment of an episode in the dynamic episodic buffer is fed to the input of the candidate episodes buffer, and the weights of its winning candidate are then copied to the static episodic buffer and replayed in the signal-encoding SOM where they generate top-down biases for sensorimotor elements. To test this ability, we exposed the trained network to 100 sequences randomly selected from the training set element by element and examined the prediction about the possible completion of the

Table 1 Transcription rules for episodes of different types

Episode :	Intransitive IntrWithPPComplement Transitive SimpleCausative CausativeWithPP
Intransitive :	Agent IntrVerb.
IntrWithPPComplement :	Agent IntrVerb2 PP.
Transitive :	Agent Target TransVerb.
SimpleCausative :	Agent Target CausativeVerb ResultVerb.
CausativeWithPP :	Agent Target CausativeVerb ResultVerb2 PP.
Agent :	AnimateObj
Target :	AnimateObj InanimateObj
PP:	Preposition Landmark
Landmark :	Target
AnimateObj :	man dog cat
InanimateObj :	cup ball chair
Preposition :	under behind near
IntrVerb :	die walk lie sneeze sit sleep smell run snore breathe
IntrVerb2 :	sneeze sit sleep smell run
TransVerb :	grab hit push shove see bite hold squeeze kick hug
CausativeVerb :	caused
ResultVerb :	break stop go
ResultVerb2 :	go hide

The colon separates the head and tail of each rule, and alternative tails are separated with |. Words starting with capital letters are non-terminal symbols. There are 35 terminal symbols, corresponding to individual sensorimotor signals/operations and an end-of-episode (.)

sequence after each element. The completion of fragments is inherently ambiguous, as there may be more than one possible continuation consistent with the episodes seen during training.

We evaluated several measures of success in prediction. The measures were separately evaluated for different relative fragments lengths and also in total (regardless of the relative fragment length).

- Grammaticality:** We evaluated the percentage of cases when the predicted sequence of SM signals represents a correctly formed episode (i.e. with a correct order of constituents).
- Compatibility:** Compatibility expresses the percentage of cases when the predicted episode starts with the presented fragment (for example, for the fragment DOG CAT..., the episode DOG CAT BITE is compatible, but DOG MAN BITE is not).
- Matches:** The simplest measure expresses the percentage of cases when the predicted episode is equal to the one being presented in fragments. However, because an episode fragment can be compatible with multiple episodes seen during training, there is an inherent ambiguity and a single “correct” answer does not have to exist.

Rank:

Because the model cannot know which of the episodes compatible with the presented fragment the fragment is part of, the best the model can do is to predict according to the frequencies of episodes seen during training. To verify this expectation, we use a theoretical variable-length Markov model (VLMM) trained on the training set. The VLMM yields possible continuations ranked by their frequencies for any fragment seen during training. Hence, we can compare the prediction of our model with that of the VLMM; if the prediction is found among those predicted by the VLMM, we record its rank.⁵ In the ideal case where the model always predicted the most frequent continuation, the average rank would be equal to 1.

⁵ If the predicted episode ranked among multiple episodes predicted by the VLMM with equal frequency, e.g. occupying 2nd–4th position, its rank would be the upper bound, i.e. 2.

Rank base: The drawback of the rank measure is that not all predicted episodes can be evaluated for rank: sometimes, the predicted episode is not found among those predicted by the VLMM, because it is either ungrammatical, or correctly formed but not found in the training set. That is why the rank must always be considered in connection with the measure we call “rank base”—the percentage of predicted episodes that were found among those predicted by the VLMM (hence could be evaluated for rank).

The results are summarised in Table 2 (top). To evaluate these results, we compared them to a baseline case of random predictions. The probability of generating a correctly formed episode by random chaining of sensorimotor elements is very low, so we used a generator that randomly selected episodes from the set of all possible *grammatical* episodes. The results are summarised in Table 2 (bottom).

We can see that the random predictor has much lower compatibility, much smaller rank base and the predictions in the rank base have much lower rank. The random predictions hardly ever match the episodes being presented. If we have a look at our model, the compatibility of prediction is very high (higher for shorter fragments, because it is easier to be compatible with one- to two-word fragments than with almost the whole sequence). In terms of grammaticality, rank, and rank base, predictions are very good for longer (50 %+) fragments: over 90 % of predictions are grammatical and the average rank is close to 1, which means the model almost always predicts the continuation most frequent in the training data. For shorter fragments, performance is not so good. It should be borne in mind that our network is performing pattern completion, rather than simply predicting the next element in a time-series; a reasonable amount of the pattern must be presented in order for it to work. In the next section, we discuss methods for improving performance on shorter fragments.

Relation to Neural Activation Data

As discussed in “[Representation of Prepared Sequences in Prefrontal Cortex](#)”, PFC stores prepared sequences in several ways. We examined the properties of representations in the trained network to see how they corresponded to representations identified in monkey PFC.

Some PFC cells encode individual operations in a prepared sequence, in a way which takes into account the sequential context they appear in (see, for example, Barone

and Joseph [26]). An example of such a cell is given in Fig. 3a. Inspecting units in the signal-encoding SOM shows that they have this property (see Fig. 4). Since the dynamic and static episodic buffers are isomorphic to the signal-encoding SOM, units in these media also have this property. Figure 3b shows responses of a particular unit in the dynamic episodic buffer for the concept *MAN* in different sequential contexts.

Some PFC units encode individual operations in a prepared sequence in a format where relative activation levels indicate the serial order in which operations will be executed [27]. Of these units, some have activity which changes dynamically during execution of a prepared sequence, being maximal before execution of the action they encode and being inhibited thereafter. Others are invariant during execution of a planned sequence. Units in the dynamic episodic buffer have the former property, and units in the static episodic buffer have the latter property. An example of PFC cells with the former property is given in Fig. 5a. An example of cells from the dynamic episodic buffer with the same property is given in Fig. 5b.

Finally, some macaque PFC neurons hold representations of candidate alternative prepared sequences, which appear to compete against one another so that a winning sequence can be selected [29]. In our model, the candidate episodes buffer holds units that have this function. Table 3 shows the five most active candidates in the candidate episodes buffer as a response to the presentation of *DOG BALL* and *DOG BALL CAUSE* fragments.⁶

Storage Capacity

Unlike in most traditional feature maps, the number of different sequences that can be stored with m units in the dynamical episodic buffer is not linear, but exponential in m , thanks to its distributed nature. The combination of the MSOM architecture of the signal-encoding SOM with the isomorphic activity gradient-based dynamical episodic buffer is functionally similar to the SARDNET architecture [41]: similar elements occupy similar positions in the map (moreover, thanks to the recurrent nature of MSOM, similar elements *in similar contexts* occupy similar positions) and each sequence is represented by a distributed pattern. Hence, the SARDNET capacity analysis [41] also applies to our model: if the maximum length of a SM sequence is l and each position can be occupied by max. k elements, the number of possible sequences is approximately k^l . Because sequences can contain repeated elements and these must be represented by different units, lk units are necessary to represent k^l

⁶ Candidates were determined by top-down reconstruction, i.e. replayed as a temporal sequence in the aggregate SM signal layer.

Table 2 Prediction performance of our model (top) and of a baseline model picking up a correctly formed episode at random (bottom), on initial fragments of different lengths of 100 episodes from the training set

Fragment length	0–25 %	25–50 %	50–75 %	75–100 %	100 %	Total
Our model						
Grammaticality	42.1 % (37.1)	74.7 % (7.1)	90.1 % (3.2)	95.4 % (1.9)	96.2 % (1.6)	81.3 % (7.0)
Compatibility	100.0 % (0.0)	98.4 % (1.2)	95.2 % (2.9)	87.2 % (5.4)	90.3 % (4.0)	94.6 % (2.2)
Matches	0.0 % (0.0)	5.0 % (1.2)	41.8 % (2.5)	84.4 % (6.8)	88.5 % (4.4)	41.3 % (2.1)
Rank	5.06 (7.04)	2.19 (0.66)	1.27 (0.10)	1.02 (0.02)	1.00 (0.00)	1.92 (1.11)
Rank base	42.0 % (37.2)	74.7 % (7.1)	86.9 % (3.1)	85.1 % (5.7)	88.5 % (4.4)	77.5 % (7.2)
Baseline random model						
Grammaticality	100.0 % (0.0)	100.0 % (0.0)	100.0 % (0.0)	100.0 % (0.0)	100.0 % (0.0)	100.0 % (0.0)
Compatibility	34.5 % (4.6)	7.2 % (1.7)	1.4 % (0.8)	0.1 % (0.3)	0.0 % (0)	7.0 % (0.8)
Matches	0.0 % (0.0)	0.0 % (0.0)	0.0 % (0.0)	0.0 % (0.0)	0.0 % (0.0)	0.0 % (0.0)
Rank	23.91 (3.43)	18.97 (7.89)	2.40 (1.74)	1.00 (0)	1.00 (0)	21.46 (2.1)
Rank base	12.7 % (5.6)	2.7 % (5.3)	0.4 % (5.0)	0.1 % (1.5)	0.1 % (0.3)	2.6 % (3.2)

Results are averaged over ten different simulation runs. Numbers in parentheses represent standard deviations

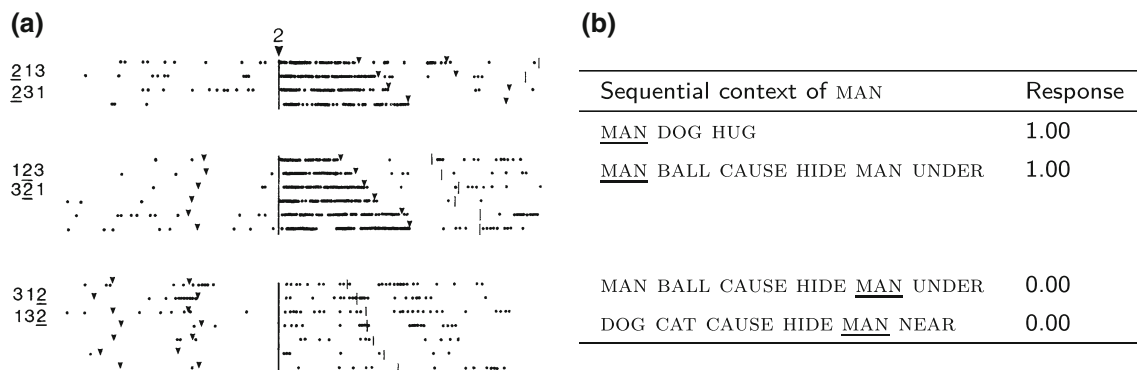


Fig. 3 a Output from a PFC cell in Barone and Joseph’s [26] experiment which fires when the monkey executes eye movement 2 in some sequential contexts (sequences [2, 1, 3], [2, 3, 1], [1, 2, 3] and [3, 2, 1]) but not others (sequences [3, 1, 2] and [1, 3, 2]). Picture reproduced from Barone P, Joseph JP. Prefrontal cortex and spatial sequencing in macaque monkey. *Exp Brain Res.* 1989;78:447–64 with kind permission from Springer Science and Business Media (Original caption: Activation of a fixation cell. All rasters are aligned with onset of the saccade towards target “2”. In the first raster, trials

corresponding to sequences in which target “2” ranked first (213, 231) are gathered, in the second raster, sequences in which it ranked second (123, 321) and in the third one, sequences in which it ranked third (312, 132). The other conventions are the same as in Fig. 3. Note that the cell shows a sustained activation whenever the animal fixates target “2”. **b** The *top-leftmost* unit of the dynamic episodic buffer fires for the signal MAN in some sequential contexts ([MAN DOG HUG], [MAN BALL CAUSE HIDE MAN UNDER]), but not others ([MAN BALL CAUSE HIDE MAN UNDER], [DOG CAT CAUSE HIDE MAN NEAR])

sequences with unique distributed patterns. We ran experiments with the same training sets as before (500 sequences built from 35 SM elements), focused on the influence of the size of the dynamical episodic buffer (which equals to the size of the isomorphic static episodic buffer and signal-encoding SOM). For our input data, the maximum sequence length is six, the number of items being sequenced is 35, so the expected size for a network with perfect storage is $6 \times 35 = 210$ units. Figure 6 shows the effect of the reduced capacity on the immediate serial recall: the size of the network where

performance comes close to 100 % is roughly 225 (15×15). The performance is equally good for unseen sequences that sample the space of possible sequences beyond the training set. Hence, the experimental results correspond to our theoretical analysis.

Regarding the capacity of the candidate episode buffer, this is linear in the number of units. However, as we discuss in ‘Coverage Issues’ below, this is not a medium which needs to represent all possible episodes, rather its role is to represent expected or desired episodes or episode types. In this role, its limited capacity works well; due to

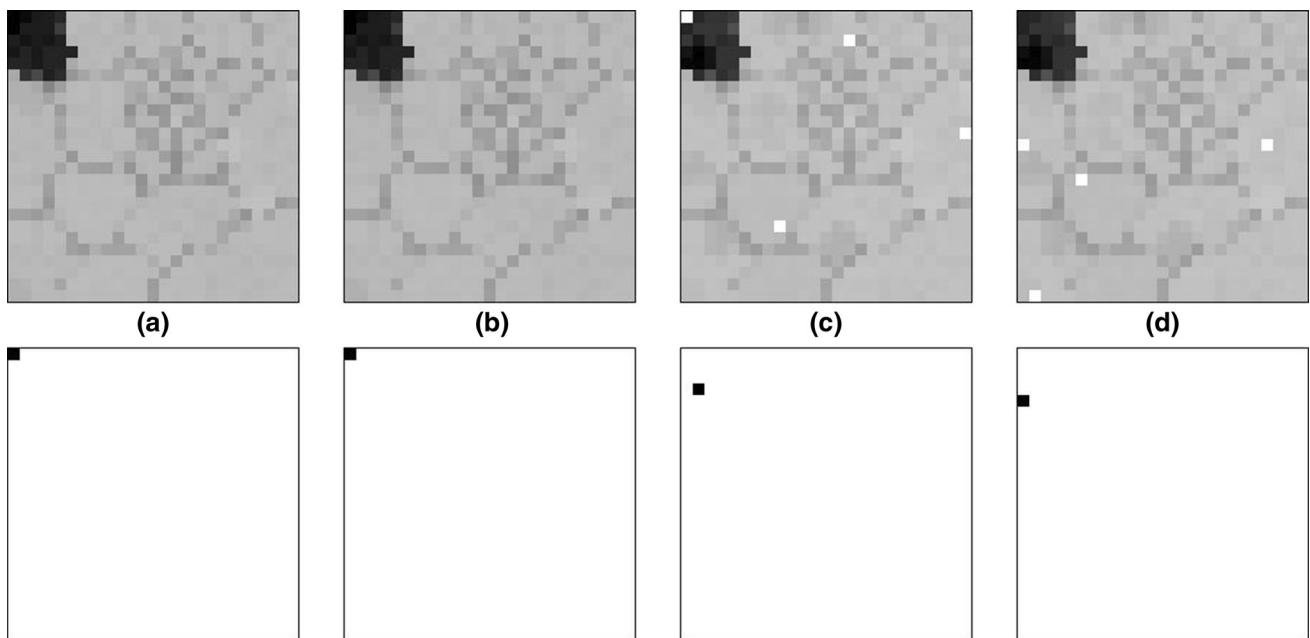


Fig. 4 Exposure to a SM signal MAN in different contexts leads to different activity profiles in the signal-encoding SOM (*top row*). A unit with the highest activity is depicted in the corresponding frame in the *bottom row*. This winning unit then activates a unit with the isomorphic position in the dynamic episodic buffer. The different

contexts for the SM signal MAN were as follows: MAN DOG HUG (**a**), MAN BALL CAUSE HIDE MAN UNDER (**b**), MAN BALL CAUSE HIDE MAN UNDER (**c**), and DOG CAT CAUSE HIDE MAN NEAR (**d**). The activities are *colour-coded* in shades of grey from *white* (activity = 0) to *black* (activity = 1)

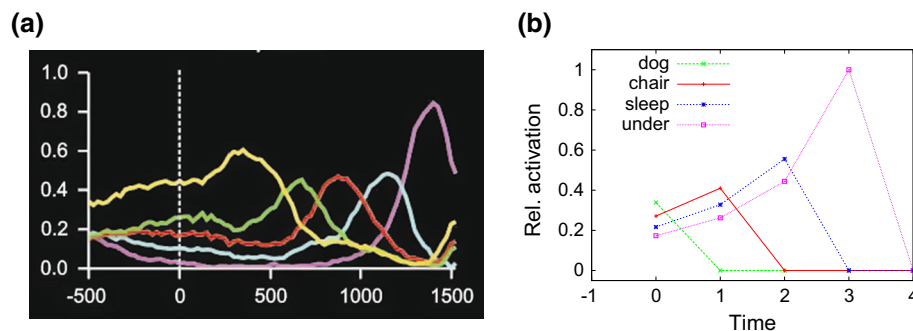


Fig. 5 a Output from four PFC cells in Averbeck et al.'s [27] experiment, which encode a prepared motor sequence, executed at time 0. The relative activation levels of the cells at time 0 indicates the order in which the four movements will be executed. Picture reproduced with permission from Averbeck B, Chafee M, Crowe D, Georgopoulos A. Parallel processing of serial movements in prefrontal cortex. PNAS. 2002;99(20):13172–7. Copyright (2002)

the properties of the SOMs on the route, if the capacity is exceeded, highly similar episodes will map to the same unit and will be conflated together.

Discussion

This paper contains two proposals. The first is a general proposal about how episodes are represented in semantic WM. We propose that episodes are stored in this medium

National Academy of Sciences, U.S.A. **b** Relative activation of four units in the dynamic episodic buffer encoding the SM sequence DOG CHAIR SLEEP UNDER (see Fig. 2), showing a similar temporal profile. The use of a relative measure of activation provides a simple approximation of the dynamics of a continuous circuit implementing winner-take-all and inhibition-of-return (see, for example, the review in [40])

as *prepared sensorimotor routines*. If this is the case, then models of WM episode representations can draw on the rich experimental tradition investigating representations of prepared behavioural sequences in prefrontal cortex. Our second proposal is a computational model of episode representations that does exactly this. The model represents episodes as prepared sensorimotor sequences; it stores these sequences in three formats that have also been found in monkey prefrontal cortex. We now assess these two proposals.

Table 3 Activity levels of the top five most active units representing episodes in the candidate episodes buffer predicted from the fragments DOG BALL (left) and DOG BALL CAUSE (right)

Activity	Reconstructed sequence
Sequence fragment: DOG BALL	
0.30	DOG BALL PUSH
0.27	DOG BALL SEE
0.27	DOG BALL GRAB
0.26	DOG BALL KICK
0.25	DOG BALL HIT
Sequence fragment: DOG BALL CAUSE	
0.33	DOG BALL CAUSE GO
0.32	DOG BALL CAUSE STOP
0.32	!DOG BALL CAUSE GO CAT BALL CAT CAUSE GO
0.29	DOG BALL CAUSE HIDE DOG NEAR
0.29	DOG BALL CAUSE HIDE MAN UNDER

The exclamation mark ‘!’ denotes an ‘ill-formed’ episode representation.

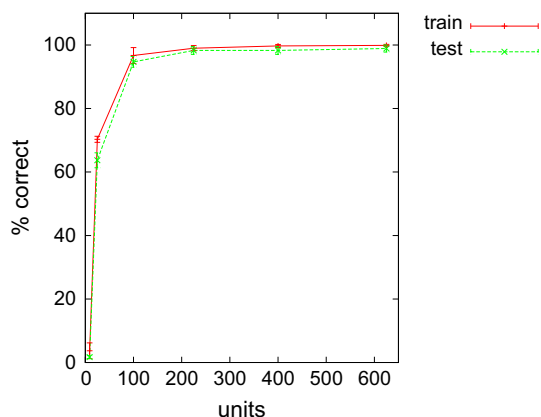


Fig. 6 Effects of the size of the dynamical episodic buffer on the success in immediate serial recall for 100 sequences from the training set (*train*) and 100 unseen sequences (*test*). We tested sizes 3×3 , 5×5 , 10×10 , 20×20 , 25×25 units. The results were averaged over three runs with different initial weights and training sets (the graph also shows standard deviations, but they are very small). We also analysed incorrectly recalled sequences: in all the cases except 3×3 , only one element of the sequence was incorrectly recalled—all the others were correct

Representing Episodes as Prepared Sequences

Coverage Issues

An obvious issue for discussion is whether a scheme that represents episodes as sequences is sophisticated enough to represent *all possible episodes*. We begin by considering this question.

We certainly do not expect that the *candidate episodes buffer* is able to represent all possible episodes. The role of this medium is just to store a distribution of possible episodes, highlighting the most likely or desirable episodes. The capacity to represent all possible episodes resides elsewhere in the network, in the dynamic episodic buffer, which stores a prepared sequence as a set of signals with different levels of activation. This activity-based sequence representation is quite productive for atomic episodes: the network can encode and reproduce 98.9 % of unseen episodes (see “[Immediate Serial Recall](#)” section).

But we have not yet considered how a sequence-based scheme could represent episodes *in which other episodes are nested*. As well as atomic episodes like *The dog barked*, we must also presumably be able to represent *The dog [which chased Mary] barked*, *Sally tells Bill that [the dog barked]*, and so on. This requirement combinatorially expands the number of episodes that the system must represent.

In response, we suggest that it is important to distinguish *episode* representations from representations of ‘sentence meanings’. A sentence meaning may involve multiple episode representations. For instance, when we entertain the meaning of *The dog [which chased Mary] barked*, this presumably involves representing two *separate* episodes: the barking episode and the chasing episode. Secondly, our general conception of semantic representations as stored sensorimotor sequences means that we do not necessarily have to model the meanings of all episodes in a nested sentence *in a single static pattern of activity*. An alternative open to us is to model complex meanings of this kind by positing structures that cause the network to activate a *sequence* of simple episode representations, in a manner that indicates what relationship they have to each other. For example, to model *The dog [which chased Mary] barked*, we can begin simply by rehearsing the matrix episode *The dog barked*. If, after having activated the initial representation of a token dog, we choose to enter a new cognitive mode in which object representations are associated with salient episodes in which they have participated, this could cause us to temporarily activate a new subordinate episode *The dog chased Mary* and rehearse this episode within the original rehearsal process, resuming the top-level rehearsal process when it is complete. (A scheme of this kind is used in Miikkulainen’s [42] treatment of relative clauses.) This device of interrupting processing is not available to schemes which represent the complete meaning of a nested proposition as a single static pattern of neural activity. We see this as a strong advantage of representing episodes as sequences. We have also used sequentially structured episode representations to provide an interesting representation of nested sentential complements such as *Sally tells*

Bill that [the dog barked]; for details, see Caza and Knott [19], Knott [20].

Sentence Processing

Any model of WM episode representations must ultimately include an account of how these representations participate in sentence processing. Baddeley [2] proposes links between the episodic buffer and the phonological buffer to support sentence processing, but does not suggest in any detail how these might be used. However, if episodes are represented in WM as prepared sequences, this opens the way for some interesting new proposals about sentence processing.

We consider first sentence generation. In most neural models of generation, the semantic inputs to the process are static patterns of activation (see, for example, [43, 44]). Since a sentence is a sequence of linguistic units, these schemes see generation as involving a process of linearisation, in which a static, declarative semantic representation is converted to a sequence of words. In our account, episode representations in WM already have sequential structure, allowing an alternative model of generation, founded on the active rehearsal of an episode-denoting sensorimotor sequence. A detailed model of sentence generation based on this idea is given in [45]. In this model, generating a sentence involves replaying a selected WM episode, in a special mode where sensorimotor signals can trigger learned phonological side effects. During this replay process, an interesting mixture of sustained and transient signals is evoked: in particular, there are tonically active representations of each action in the planned sequence in the static episodic buffer throughout the replay process. These tonic representations permit a neat account of the extended syntactic domain of verbs. Verbs can appear at various different positions in the structure of a clause, and they can carry inflections signalling agreement with arguments at distant positions in the clause (for instance, subjects). The neural basis for this non-locality is currently a complete mystery. But if sentences are produced by replaying a prepared sensorimotor routine and if verbs and their inflections are produced from planned motor and attentional action representations which are tonically active during replay, we have a promising explanation of this non-locality: the semantic representations from which inflected verbs are generated are active throughout the generation process and can be produced at any time.

The WM episode network also has interesting uses in models of sentence interpretation. Interpreting a sentence involves selecting a single meaning from a large set of alternative possible meanings, a process in which probabilistic calculations play an important role [46]. In our network, the candidate episodes buffer represents a set of alternative

episode representations from which a winner can be chosen. This layer could be a useful medium for representing a distribution of possible sentence interpretations, reflecting the structure of the incoming sentence as it arrives incrementally, but also the hearer's general experience of episodes, and generating expectations about the rest of the sentence. An attractive feature of the candidate episodes buffer is that it allows a single winning episode to be selected at any point during experience of an episode, which can generate expectations about the remainder of the episode. If these expectations are not forthcoming, it also permits the inhibition of the predicted episode and the selection of an alternative candidate—a process that could perhaps model the 'backtracking' operation that occurs during the interpretation of a garden-path sentence. This is an idea we are pursuing in current work.

Predictions About Dysfunction

Our model sees a WM episode representation as a high-level action plan implemented in dlPFC—a plan to reproduce the sequence of attentional and motor operations through which a given episode was experienced. If there is dysfunction in the prefrontal system that plans sequences of attentional and/or motor operations, we therefore predict two other types of impairment. First, we predict impaired sentence processing abilities, since in our model, the WM episode buffer is the medium that holds the semantic representation of sentences. There are several studies indicating that impairments in attentional sequencing tasks correlate with impairments in sentence processing. For instance, Dispaldro et al. [47] show that children with specific language impairment also have impaired performance on visual attentional sequencing tasks. Second, we predict impaired storage of episodes in WM and LTM, since in our model, the process of creating a LTM representation of an episode relies on it first being stored in the WM episodes buffer. There is good evidence that transcranial magnetic stimulation over the dorsolateral PFC disrupts both WM representations and the encoding of representations in episodic LTM (see Balconi [48]). In summary, studies of dysfunction support the idea that the system involved in preparing sequences of attentional and motor actions also has a role in encoding WM representations that interface to language and to LTM.

The SOM Model of Episode Representations

In this section, we assess the neural network model for episode representations more technically and compare it to existing models.

The Use of Localist Representations in the Network

One issue to discuss is the network's use of localist representations. Both input signals and episodes are represented in a localist scheme. This is clearly unrealistic as a model of biological neural networks; in addition, the network may be better at representing generalisations over episodes if it works with distributed input representations. To investigate how our architecture works with distributed representations, we implemented a revised version of the network using distributed inputs, and distributed encodings of signals and episodes in the SOM layers [37]. The network correctly recalled 94.6 % (SD 3.8 %) of seen sequences, 93.3 % (SD 3.7 %) of unseen sequences, and 93.3 % (SD 6.3 %) of sequences containing repetitions, when presented with the same sets of sequences as the localist model, a drop of around 5 %. However, the distributed model's performance in predicting episodes is slightly better⁷ than that of the localist model. While there is still room for improvement, these experiments provide some indication that our proposed architecture can be adapted to operate with distributed representations.

Comparison with Other Models of Sequence Representation

There are numerous network models of WM for sequences. Some represent the order of elements by their associations with a context signal that evolves independently as a function of time [49–51] or a context signal that evolves as a function of its previous state [52]. Others represent the order of elements by their level of activation [53] or in sustained patterns of activation within a recurrent neural network [54]. These models share several features with ours: for instance, the use of a recurrent neural network for context-dependent representation of sequence elements in the signal-encoding SOM, and activation gradient and inhibition of the recent winner in the dynamic episodic buffer. However, the cited models are explicitly stated as models of *phonological* WM. We follow Baddeley [2] in distinguishing between phonological WM and WM for episodes. Most importantly, this means our model does not have to reproduce the classic effects found in immediate recall of phonological sequences, such as the primacy and recency effects, list-length effect, transposition gradient. Empirically, our focus is on modelling the neural sequence preparation mechanisms found in monkeys, which it does quite successfully. The sequences we use represent certain internal thematic role structures of episodes (see

“**Training**” section); hence, the order of elements is far from arbitrary, and the classical immediate serial recall behavioural data are not relevant.

There are some computational models which propose the same mechanism for both phonological WM and prepared action sequences—see, in particular, Rhodes et al. [28]. We certainly envisage similarities between the mechanisms subserving these tasks. (In particular, they both appear to involve competitive queuing.) But our suggestion is that they are separate, although, as Baddeley suggests, there are links between them, which support sentence processing.

Perhaps the model most closely related to ours is that of Dominey and colleagues. Dominey et al. [55] developed a neural network model of sequence preparation which is explicitly based on data from single-cell recordings in monkeys by Barone and Joseph [26]. This model was later recruited in a neural network for sentence processing, modelling how sequences of words are mapped onto working memory episode representations [56]. However, while this model links prefrontal sequence preparation mechanisms with working memory episode representations, it does not construe working memory episode representations *themselves* as prepared sensorimotor sequences. In addition, the model does not incorporate the interesting prefrontal sequence preparation mechanisms identified by Averbeck et al. [27]. Nonetheless, there are perhaps ways of integrating our model with that of Dominey et al. In their account of sentence interpretation, nothing hangs on the particular format of episode representations they choose, so there is perhaps scope for combining their account of sentence interpretation with our account of episode representations.

Comparison with Other Models of Episode Representation

There are many neural network models of episode representation. A key question for these models is how to associate representations of participants in an episode with the roles they play (e.g. ‘agent’ and ‘patient’). This binding has been implemented using many techniques, including synchronised neural firing (e.g. [57]), potentiated synaptic associations (e.g. [43, 58]), high-dimensional representations (e.g. [59]), and self-organising maps (e.g. [7]). One unifying aspect of all these models is that they use explicit representations of semantic roles: roles like ‘agent’ and ‘patient’ are encoded as patterns of activity, in basically the same way as the items they are bound to. It is not clear whether humans represent roles explicitly in this way. If they do, these explicit representations are certainly enormously active in the cognitive system—common roles like ‘agent’ and ‘patient’ feature in nearly every episode. We would certainly predict specific types of failure if one of these representations were to be damaged. A key difference

⁷ The distributed model performed better by 3.7 % in total grammaticality, 0.2 % in compatibility, 0.8 % in matches, 0.09 in rank, and 5.1 % in rank base [37].

in our model is that semantic roles are represented implicitly, by particular positions within a sensorimotor sequence. This encoding avoids a troublesome assumption with existing proposals. Instead, the notions of ‘agent’, ‘patient’, and so on, are defined through direct reference to the processes by which episodes are actually experienced.

Another model we want to mention is the syntagmatic paradigmatic model of Dennis [60]. Although stated as a memory-based model of verbal processing, it deals with tasks that are similar to those implemented by our model. Dennis’ model represents propositional and thematic role information in an interesting way. The model consists of a long-term memory for sequences (sentences from a corpus) and paradigmatic relations (information about which words appear in similar contexts). A set of words that appear in a similar context implicitly represents a specific thematic role (such as *lover* or *lovee*). When prompted with a new sequence or query in working memory, the system uses Bayesian inference to retrieve the most similar sentences from the long-term memory together with probabilities for role bindings/substitutions based on paradigmatic relations.

The model of episode representation that comes closest to ours is perhaps that of Reynolds et al. [61]. This model is also founded on the assumption that experiencing an episode is a process with a well-defined sequential structure. However, in their model, this sequential structure is modelled as a continuous trajectory in a high-dimensional space of relatively low-level perceptual features; by contrast, our model envisages a sequence of more discrete, high-level operations. We do not see these models as inconsistent; however, they can be thought of as describing processing at two different levels of representation. It is an interesting question whether they can be combined in some way.

Conclusion

In this paper, we presented a model of the storage of episode representations in working memory. The main theoretical novelty is a proposal that episodes are stored as prepared sensorimotor routines. Our proposed model also implements a novel method for representing multiple competing sequences in the same medium at the same time. Naturally, the current implementation must be refined in many ways before it can serve as a proper model of neural mechanisms; for instance, it should use distributed representations, and it should be adapted to work with spiking neurons. And while the model is already aligned with data about sequence preparation in monkeys, its predictions about equivalent data in humans should also clearly be examined. These are interesting avenues for further research.

Acknowledgments This research was supported by NZ Marsden Fund, and partially supported by grants VEGA 1/0898/14 and KEGA 076UK-4/2013 for Martin Takac. We are grateful to Lubica Benuskova and Igor Farkas for helpful discussions.

References

1. Baddeley A, Hitch G. Working memory. In: Bower G, editor. The psychology of learning and motivation: advances in research and theory (Vol. 8). New York: Academic Press; 1974. p. 47–90.
2. Baddeley A. The episodic buffer: a new component of working memory? *TICS*. 2000;4(11):417–23.
3. Tulving E, Markowitsch H. Episodic and declarative memory: role of the hippocampus. *Hippocampus*. 1998;8:198–204.
4. Abraham W, Logan B, Greenwood J, Dragunow M. Induction and experience-dependent consolidation of stable long-term potentiation lasting months in the hippocampus. *J Neurosci*. 2002;22:9626–34.
5. Levelt W. *Speaking: from intention to articulation*. Cambridge: MIT Press; 1989.
6. Haarmann H, Cameron K, Ruchkin D. Short-term semantic retention during on-line sentence comprehension. Brain potential evidence from filler-gap constructions. *Cognit Brain Res*. 2003;15:178–90.
7. Mayberry M, Miikkulainen R. Incremental nonmonotonic sentence interpretation through semantic self-organization. The University of Texas at Austin: Department of Computer Sciences; 2008; AI08–12.
8. Ballard D, Hayhoe M, Pook P, Rao R. Deictic codes for the embodiment of cognition. *Behav Brain Sci*. 1997;20(4):723–67.
9. Knott A. *Sensorimotor cognition and natural language syntax*. Cambridge: MIT Press; 2012.
10. Knott A. Syntactic structures as descriptions of sensorimotor processes. *Biolinguistics*. 2014;8:1–52.
11. Wallenstein G, Eichenbaum H, Hasselmo M. The hippocampus as an associator of discontinuous events. *Trends Neurosci*. 1998;21:317–23.
12. Lee A, Wilson M. Memory of sequential experience in the hippocampus during slow wave sleep. *Neuron*. 2002;36:1183–94.
13. Cutsuridis V, Hasselmo M. Spatial memory sequence encoding and replay during modeled theta and ripple oscillations. *Cognit Comput*. 2011;3:554–74.
14. Haggard P. Human volition: towards a neuroscience of will. *Nat Rev Neurosci*. 2008;9:934–46.
15. Tessitore G, Prevede R, Catanzariti E, Tamburrini G. From motor to sensory processing in mirror neuron computational modelling. *Biol Cybern*. 2010;103:471–85.
16. Chambon V, Wenke D, Fleming S, Prinz W, Haggard P. An online neural substrate for sense of agency. *Cereb Cortex*. 2012;23:1031–7.
17. Tipper S, Lortie C, Baylis G. Selective reaching: evidence for action-centered attention. *J Exp Psychol Hum Percept Perform*. 1992;18:891–905.
18. Johansson R, Westling G, Backstrom A, Flanagan J. Eye-hand coordination in object manipulation. *J Neurosci*. 2001;21(17):6917–32.
19. Caza G, Knott A. Pragmatic bootstrapping: a neural network model of vocabulary acquisition. *Lang Learn Dev*. 2012;8:1–23.
20. Knott A. How infants learn word meanings and propositional attitudes: a neural network model. In: Hung TW, editor. *Lang action*. Berlin: Springer; 2014. p. 107–24.
21. Lee-Hand J, Knott A. A model of causative actions: from motor learning to syntactic structure. In: *Proceedings of the 35th annual*

- meeting of the cognitive science society. Berlin; 2013; p. 2849–2854.
22. Curtis C, D'Esposito M. Persistent activity in the prefrontal cortex during working memory. *Trends Cognit Sci.* 2003;7(9):415–23.
 23. Barsalou L. Grounded cognition. *Ann Rev Psychol.* 2008;59:617–45.
 24. Glenberg A, Gallese V. Action-based language: a theory of language acquisition, comprehension, and production. *Cortex.* 2012;48:905–22.
 25. Tulving E. Elements of episodic memory. New York: Oxford University Press; 1983.
 26. Barone P, Joseph JP. Prefrontal cortex and spatial sequencing in macaque monkey. *Exp Brain Res.* 1989;78:447–64.
 27. Averbeck B, Chafee M, Crowe D, Georgopoulos A. Parallel processing of serial movements in prefrontal cortex. *PNAS.* 2002;99(20):13172–7.
 28. Rhodes B, Bullock D, Verwey W, Averbeck B, Page M. Learning and production of movement sequences: behavioral, neurophysiological, and modeling perspectives. *Hum Mov Sci.* 2004;23:699–746.
 29. Averbeck B, Sohn J, Lee D. Activity in prefrontal cortex during dynamic selection of action sequences. *Nat Neurosci.* 2006;9(2):276–82.
 30. Shallice T, Burgess P. Deficits in strategy application following frontal lobe damage in man. *Brain.* 1991;114:727–41.
 31. Kaller C, Heinze K, Frenkel A, Lippchen C, Unterrainer J, Weiller C, et al. Differential impact of continuous theta-burst stimulation over left and right DLPFC on planning. *Hum Brain Mapp.* 2013;34:36–51.
 32. Ziemke T. What's that thing called embodiment? In: Alterman R, Kirsh D, editors. Proceedings of the 25th Annual Conference of the Cognitive Science Society. Mahwah, NJ: Lawrence Erlbaum; 2003. p. 1134–9.
 33. Kohonen T. Self-organized formation of topologically correct feature maps. *Biol Cybern.* 1982;43:59–69.
 34. Kahneman D, Treisman A, Gibbs B. The reviewing of object files: object-specific integration of information. *Cognit Psychol.* 1992;24:175–219.
 35. Hasegawa R, Matsumoto M, Mikami A. Search target selection in monkey prefrontal cortex. *J Neurophysiol.* 2000;84:1692–6.
 36. Fix J, Rougier N, Alexandre F. A dynamic neural field approach to the covert and overt deployment of spatial attention. *Cognit Comput.* 2011;3:279–93.
 37. Takac M, Knott A. A revised neural network model of episode representations in working memory. Dept of Computer Science, University of Otago; 2014; OUCS-2014-03.
 38. Strickert M, Hammer B. Merge SOM for temporal data. *Neurocomputing.* 2005;64:39–71.
 39. Mayr U, Keele S. Changing internal constraints on action: the role of backward inhibition. *J Exp Psychol Gen.* 2000;129(1): 4–26.
 40. Vickers D, Lee MD. Dynamic models of simple judgments: II. Properties of a self-organizing PAGAN (parallel, adaptive, generalized accumulator network) model for multi-choice tasks. *Nonlinear Dyn Psychol Life Scie.* 2000;4(1):1–31.
 41. James DL, Miikkulainen R. SARDNET: a self-organizing feature map for sequences. In: Tesauro G, Touretzky DS, Leen TK, editors. Advances in neural information processing systems 7. Cambridge: MIT Press; 1995. p. 577–84.
 42. Miikkulainen R. Subsymbolic case-role analysis of sentences with embedded clauses. *Cognit Sci.* 1996;20:47–73.
 43. Chang F. Symbolically speaking: a connectionist model of sentence production. *Cognit Sci.* 2002;26:609–51.
 44. Rohde D. A connectionist model of sentence comprehension and production. School of Computer Science, Carnegie Mellon University; 2002.
 45. Takac M, Benuskova L, Knott A. Mapping sensorimotor sequences to word sequences: a connectionist model of language acquisition and sentence generation. *Cognition.* 2012;125: 288–308.
 46. Collins M. A new statistical parser based on bigram lexical dependencies. In: Proceedings of the 34th meeting of the ACL. Santa Cruz; 1996; p. 184–191.
 47. Dispaldro M, Leonard L, Corradi N, Ruffino M, Bronte T, Fa-coetti A. Visual attentional engagement deficits in children with specific language impairment and their role in real-time language processing. *Cortex.* 2013;49:2126–39.
 48. Balconi M. Dorsolateral prefrontal cortex, working memory and episodic memory processes: insight through transcranial magnetic stimulation techniques. *Neurosci Bull.* 2013;29(3):381–9.
 49. Burgess N, Hitch G. Memory for serial order: a network model of the phonological loop and its timing. *Psychol Rev.* 1999;106: 551–81.
 50. Burgess N, Hitch G. A revised model of short-term memory and long-term learning of verbal sequences. *J Mem Lang.* 2006;55(4):627–52.
 51. Brown GDA, Hulme C, Preece T. Oscillator-based memory for serial order. *Psychol Rev.* 2000;107:127–81.
 52. Howard M, Kahana M. A distributed representation of temporal context. *J Math Psychol.* 2002;46:269–99.
 53. Page M, Norris D. A model linking immediate serial recall, the Hebb repetition effect and the learning of phonological word forms. *Philos Trans R Soc B.* 2009;364:3737–53.
 54. Botvinick MM, Plaut DC. Short-term memory for serial order: a recurrent neural network model. *Psychol Rev.* 2006;113(2): 201–33.
 55. Dominey P, Arbib M, Joseph JP. A model of corticostriatal plasticity for learning associations and sequences. *J Cognit Neurosci.* 1995;7(3):311–36.
 56. Dominey P, Hoen M, Inui T. A neurolinguistic model of grammatical construction processing. *J Cognit Neurosci.* 2006;18(12): 2088–107.
 57. Shastri L. Episodic memory trace formation in the hippocampal system: a model of cortico-hippocampal interaction. Berkeley: International Computer Science Institute (ICSI), UC Berkeley; 2001; TR-01-004.
 58. van der Velde F, de Kamps M. Neural blackboard architectures of combinatorial structures in cognition. *Behav Brain Sci.* 2006;29:37–108.
 59. Stewart T, Eliasmith C. Compositionality and biologically plausible models. In: Werning M, Hinzen W, editors. The Oxford handbook of compositionality. New York: Oxford University Press; 2012.
 60. Dennis S. A memory-based theory of verbal cognition. *Cognit Sci.* 2005;29(2):145–93.
 61. Reynolds J, Zacks J, Braver T. A computational model of event segmentation from perceptual prediction. *Cognit Sci.* 2007;31: 613–43.